

1 "Methods and systems for searching and associating information resources such as
2 web pages"

3
4 The present invention relates in a general manner to methods and systems for
5 managing resources such as web pages accessible via the Internet, or any other
6 types of documents, aimed on the one hand at improving the obtaining of resources
7 that are "close" to given resources, in terms in particular of centers of interest for
8 the user, and aimed on the other hand at allowing the user, in a particularly simple
9 and intuitive manner, to effect associations between resources himself, especially
10 so as to benefit therefrom during the obtaining of close resources.

11 State of the Art

12 The quantity of information potentially relevant for each individual is becoming
13 such that the present procedures for storing and searching for information are
14 scarcely adequate. Alongside systems making it possible to retrieve information
15 organized explicitly (such as "favorites") or by key words (via a search engine), it
16 would be desirable to have available a method which spontaneously proposes
17 context dependent relevant information.

18 Systems which provide relevant links (or rather "related links" to use the jargon)
19 with respect to a current page visited on the web are known. Typically these
20 systems comprise an extension to the Internet browser which communicates with a
21 remote server which provides the relevant links as a function of the current page
22 presented in the browser's main window. Typically these links are presented, in the
23 form of a list of URLs, in a window adjacent to the browser's main window.

24 However, such systems are not extended to serve as associative memory.

25 Summary of the Invention

26 An object of the present invention is to propose computer methods and systems for
27 searching for resources (especially web pages, diverse computer documents) that
28 are "close" to given resources (this notion of closeness being made explicit later),
29 and methods for the associative management of resources.

30 In particular, the invention is aimed at characterizing information elements with
31 respect to new pages which appear on the web, thus opening up the way to
32 multiple new applications of dynamic management of content with respect to the
33 user's browsing context.

34 More precisely, it is the aim of the invention that each information element be
35 associated with links on relevant web pages which characterize it and which are

1 automatically maintained up to date. It is thus possible to characterize nontextual
2 information, such as photos, sounds and animations (in flash, etc.) and dynamically
3 select the elements to be presented to the user as a function of the context of his
4 browsing which is also characterized by sets of relevant web pages. This approach
5 is suitable especially, but not exclusively, for magazines in the art of living,
6 fashion and in all other areas of "taste" where it is difficult to characterize through
7 key words the interest shown by the subscriber in an item of information (when for
8 example it represents a piece of music, a piece of art, a culinary dish, etc).

9 Another object of the invention is to associate other targeted elements, such as
10 targeted advertisements, with information elements, in exchange for an innovative
11 associative memory service offered to surfers.

12 In particular, the aim is that, typically by means of an extension of their browser
13 (extension downloadable from a given website), users can use the information
14 elements of this site as "associative memory". Thus, during the user's browsing, the
15 most relevant element of the site with respect to the web page visited - as well as
16 with respect to the browsing context - will be presented to him spontaneously; the
17 user will then be able to drag and drop onto this element any resource from his
18 computer, such as the icon of a file of the client station, or else the URL of a web
19 page, so as to store it. Thereafter, each time he visits any web page which is
20 relevant with respect to this element, the resource that he had stored will be
21 presented to him spontaneously, together with the resources (such as
22 advertisements) that the author of the element had himself associated with the
23 element. The advertisements presented will thus correspond to the current centers
24 of interest of the user and are provided in exchange for a new associative memory
25 service.

26 The invention is aimed moreover at harnessing modern user interfaces to create, in
27 a particularly simple and intuitive manner, associations between information
28 resources (web pages, or document files) especially within the framework of the
29 above objectives.

30 The invention proposes according to a first aspect a method for determining
31 relevant additional resources with respect to a given set of starting resources,
32 characterized in that it comprises the following steps:

- 33 a) identifying a set of citing resources that consist of all the resources having a
34 link to at least one of the starting resources,
- 35 b) forming a set of candidate resources that consists of the set of resources
36 cited by the citing resources,
- 37 c) for each candidate resource, calculating a candidate resource relevance

1 score between said candidate resource and the set of starting resources on the basis
2 of the existence of links situated in the citing resources and directed toward the
3 candidate resource and toward the starting resources, and on the basis also of citing
4 resource relevance scores assigned to each of the citing resources,

5 d) for each citing resource, recalculating a citing resource relevance score on
6 the basis of the existence, in the citing resource in question, of links to the
7 candidate resources and on the basis also of the candidate resource relevance
8 scores allocated to the candidate resources in step c),

9 e) repeating as appropriate step c) and step d) as appropriate one or more
10 times followed by step c),

11 f) determining said relevant additional resources as being the candidate
12 resources which exhibit the best candidate resource relevance scores (and as
13 appropriate also the citing resources which exhibit the best citing resource
14 relevance scores).

15 The relevance score calculation performed in step c) comprises the calculation of a
16 plurality of sums of citing resource relevance scores, each sum advantageously
17 comprising only the relevance scores of the citing resources comprising a link to a
18 given resource consisting of the candidate resource or a starting resource.

19 In a preferred manner, the above method also comprises the calculation of at least
20 one sum of citing resource relevance scores, each sum comprising only the
21 relevance scores of the citing resources comprising a link to one among a set of at
22 least two given resources, this set comprising the candidate resource and at least
23 one starting resource.

24 According to a second aspect, the invention proposes a method for determining
25 relevant additional resources with respect to a given set of starting resources,
26 characterized in that it comprises the following steps:

27 a) identifying a set of cited resources that consist of all the resources having a
28 link to at least one of the starting resources,

29 b) forming a set of candidate resources that consists of the set of resources
30 citing the cited resources,

31 c) for each candidate resource, calculating a candidate resource relevance
32 score between said candidate resource and the set of starting resources on the basis
33 of the existence of links situated in the candidate resource and in the starting
34 resources and directed toward the cited resources, and on the basis also of cited
35 resource relevance scores assigned to each of the cited resources,

36 d) for each cited resource, recalculating a cited resource relevance score on
37 the basis of the existence, in the cited resource in question, of links to the candidate

1 resources and on the basis also of the candidate resource relevance scores allocated
2 to the candidate resources in step c),
3 e) repeating as appropriate step c) and step d) as appropriate one or more
4 times followed by step c),
5 f) determining said relevant additional resources as being the candidate
6 resources which exhibit the best candidate resource relevance scores (and as
7 appropriate also the cited resources which exhibit the best cited resource relevance
8 scores).

9 The invention furthermore proposes a system for browsing among information
10 resources, each resource comprising at least one link activatable in a first mode by
11 an input device so as to bring about access to another information resource
12 designated by a resource identifier associated with this link, characterized in that at
13 least certain resources comprise at least one link activatable in a second mode with
14 the aid of an input device so as to send to an engine for searching for new
15 information resources a search query containing the resource identifier associated
16 with the link in question.

17 This system exhibits the following preferred but optional aspects:

18 * the input device is able to activate the link simultaneously in the first and second
19 modes.

20 * the activation of the link in the second mode is able to bring about the displaying
21 of a pre-existing query, to which the resource identifier associated with the link in
22 question is able to be added.

23 * the activation of the link in the second mode is able to display, in addition to the
24 pre-existing query, the information resource designated by said resource identifier.

25 The invention also proposes a system for searching for new information resources
26 on the basis of existing information resources, characterized in that it comprises a
27 search engine based on the analysis of links between the various resources and
28 accepting as input a query comprising a series of resource identifiers, a means of
29 selecting identifiers which is able to store a set of identifiers (URI) of resources
30 selected one after the other by a user, and a user activatable query generating
31 means for devising a query containing the set of identifiers previously selected
32 destined for the search engine.

33 In a preferred but nonlimiting manner, the means of selection is able to store the
34 identifiers selected in a remanent manner, in such a way that the means of selection
35 can be implemented in a manner staggered over time with a view to the generation
36 of one and the same query.

1 The invention moreover proposes a method of searching for new information
2 resources on the basis of existing information resources, characterized in that it
3 comprises the implementation of a search engine based on the analysis of links
4 between various resources and accepting as input a query comprising a series of
5 resource identifiers and in that it comprises the following steps:

- 6 - selection of identifiers (URI) of resources one after the other by a user;
- 7 - generation of a query containing the set of identifiers previously selected
8 destined for the search engine.

9 There is also proposed a method of searching for new information resources on the
10 basis of existing information resources, characterized in that it comprises the
11 implementation of a search engine based on the analysis of links between various
12 resources and accepting as input a query comprising a series of resource identifiers
13 and in that it comprises the following steps:

- 14 - generation of a query containing a set of identifiers of resources previously
15 stored in one and the same group of resource identifiers individual to a user,
16 destined for the search engine,
- 17 - generation of a signaling for the attention of the user when at least one new
18 resource identifier belonging to the group in question has been found by the
19 engine.

20 According to a preferred aspect of the above method, each group of resource
21 identifiers is represented by a graphical object on a display device of the user, and
22 in that said signaling is carried out at least by change of appearance of this
23 graphical object.

24 The invention furthermore proposes a method of managing resources in a
25 computer system provided with a display screen and with an input device for
26 cursor movement and actuation such as a mouse, each resource possessing a
27 representation displayed on the screen in such a way as to be able to be moved
28 with the aid of the input device, method characterized in that it comprises the
29 following steps:

- 30 - movement of the representation of a first resource so as to bring it above
31 the representation of a second resource,
- 32 - followed by storage, in an associative memory for managing resources, of
33 information of association between the first and second resources.

34 Certain preferred, but optional, aspects of this method are the following:

- 35 * the movement step is performed by a drag and drop technique.
- 36 * the method furthermore comprises, subsequent to the identification of a given

1 resource in a resource consultation process, the following steps:

2 - reading of the associative memory for managing resources to determine
3 whether other resources are associated with said given resource, and

4 - if so, signaling on the display screen of the existence of the associated
5 resource or resources.

6 * the resources comprise files.

7 * the resources comprise resources accessible via a network such as the Internet.

8 * the identification of a given resource is obtained via a process for identifying
9 similar or relevant resources with respect to at least one starting resource.

10 * in the case where the reading of the associative management memory determines
11 the existence of several associated resources, the signaling step comprises the
12 ordered signaling of at least part of said several associated resources.

13 * the ordered signaling is based on the determination of relevance scores of said
14 associated resources.

15 * the associative memory for managing resources is contained in a server
16 accessible from a plurality of individual stations in which the movement step can
17 be implemented.

18 * the associations between resources are stored user by user.

19 * the associations between resources are stored in a mutualized manner between
20 several users.

21 The invention also proposes a method for identifying on the basis of a text
22 resource, part of said resource able to constitute a pertinent query for a search
23 engine, characterized in that it comprises the following steps:

24 - removing the nonpertinent words from the text;

25 - establishing and completing a memory of links between parts of said text,
26 where a part is linked to another when it contains at least one pertinent word in
27 common;

28 - implementing a method of determining resource scores by analysis of a
29 graph of resource nodes connected by links, where each resource used in this
30 method consists of a part of the text, on the parts of the text that are thus
31 interconnected;

32 - using at least one of the text parts consisting of the candidate resources
33 determined by said method as query text or as basis for a query text.

1 Advantageously, the step of implementing the method for distilling resources is
2 performed only with text parts selected as prevalent, where the citing text parts are
3 the text parts which comprise at least one word in common with the prevalent text
4 part or parts, where a link is created from each citing text part to the prevalent text
5 part or parts, where the text parts containing at least one word also contained in the
6 citing text parts are identified, so as to form a group of co-cited text parts, and
7 where a link is temporarily created from each citing text part to each co-cited text
8 part with which said citing text part possesses at least one word in common.

9 The text parts are typically phrases.

10 According to another aspect, the invention proposes a method of managing
11 information resources such as web pages in a computer system comprising a user
12 station furnished with a display screen, each resource possessing an identifier
13 (URI) allowing its access from the user station, method characterized in that it
14 comprises the following steps:

- 15 a) declaration by the user of an association between two resources, by
16 associating with a second resource the identifier of a first resource;
- 17 b) identification of other relevant resources with respect to the second
18 resource; and
- 19 c) during access to one of the other resources (*current page*), signaling of the
20 existence of the first resource.

21 According to certain preferred but nonlimiting aspects:

22 * step b) comprises the selection of other resources that are most relevant for the
23 implementation of step c).

24 * step a) is implemented for a plurality of second resources belonging to a group,
25 and in that step b) comprises the identification of other relevant resources with
26 respect to the set of second resources of the group.

27 * step b) is triggered by the carrying out of step (a).

28 * step (b) is implemented subsequently to the access envisaged in step (c) to
29 determine whether the other resource which it has accessed is another relevant
30 resource with respect to the second resource.

31 * step (b) is implemented by supplying an identifier of the second resource to a
32 server for determining relevant resources.

33 * step (b) is implemented by identifying other relevant resources with respect to at
34 least one intermediate resource (*spot*) with respect to which the second resource is
35 predetermined as being relevant.

1 * the method furthermore comprises the displaying, in the vicinity of an area for
2 displaying resources, of representations of links to at least certain among the first
3 resources, the intermediate resources, and relevant resources with respect to the
4 intermediate resources.

5 * step (a) is implemented by acting with the aid of an input device on graphical
6 objects representative of the first and second resources.

7 The invention moreover proposes a method for identifying information resources
8 accessible via recent links (such as web pages), relevant with respect to at least one
9 given resource, characterized in that it comprises the following steps:

10 - applying a query comprising an identifier of said given resource to a
11 system for determining relevance between resources,

12 - selecting a first set of resources that are the most relevant (e.g. *best hub*
13 *scores*) with respect to said given resource,

14 - searching, through each of the most relevant resources, for the regions
15 possessing links to other resources of averagely high relevance, so-called relevant
16 regions,

17 - monitoring the appearance, in said relevant regions, of new links which
18 point to resources which were not yet known to the system, so-called new
19 resources,

20 - selecting a second set of resources having a high relevance (e.g. *best*
21 *hypertext authority scores*) with respect to said given resource,

22 - selecting the new resources which have a highest similarity of content with
23 respect to the resources of said second set of resources and according the new
24 resources selected a relevance level (*similarity authority score*) dependent on time
25 as a function of said similarity of content.

26 According to yet another aspect, the invention proposes a method for allowing
27 access by a user to relevant information entities from a starting information entity,
28 each information entity being accessible via an identifier (URI), characterized in
29 that it comprises the following steps:

30 a) providing at least one similar information entity, exhibiting a content
31 similar to that of the starting entity, and determining the identifier of the or of each
32 similar information entity, and

33 b) determining on the basis of the or each similar information entity identifier
34 a set of one or more identifiers of information entities relevant with respect to the
35 or each similar information entity.

36 Preferred, but nonlimiting aspects of the above method are as follows:

1 * the method furthermore comprises the following step:
2 c) allowing the user to access at least certain relevant information from their
3 respective identifiers.

4 * the method furthermore comprises the following step:
5 d) on the basis of the relevant information entity identifiers and of a given set
6 of extra information entities, selecting the extra entities that are most similar to the
7 relevant information entities.

8 * the method comprises an extra step of sorting the relevant information entities by
9 degree of relevance.

10 * the sorting step is preceded by a step of calculating a relevance score with
11 respect to the or each similar information entity for each of the relevant
12 information entities.

13 * each information entity consists of a page fragment written in a standardized
14 mark-up language, or of such a page as a whole.

15 * each identifier consists of a uniform resource identifier (URI) of the fragment or
16 of the page.

17 * step a) is carried out by selection by the user of one or more information entities
18 similar to the starting information entity.

19 * step a) is carried out by implementing a process for automatically determining
20 similar information entities.

21 * step a) is carried out by implementing a process for automatically determining
22 similar information entities, followed by a selection by the user of one or more
23 similar information entities from among the similar information entities
24 determined by said process.

25 * step b) is carried out by implementing a process for automatically determining
26 relevant information entities.

27 * the process for automatically determining relevant information entities comprises
28 the analysis of a graph structure of identifiers that consists of the identifiers of
29 information entities and of the identifiers designated by user activatable links
30 contained in said information entities.

31 According to another aspect of the invention, a method for determining relevance
32 scores of text units such as phrases in a textual document, comprises the following
33 steps:

34 - decomposition of the document into a plurality of text units,

- 1 - selection of at least one relevant text unit and of candidate text
- 2 units,
- 3 - determination of the set of pertinent words contained in the relevant
- 4 text unit (or units) and in each of the candidate text units,
- 5 - for each pertinent word contained in the relevant text unit (or units),
- 6 identification of the candidate text units citing this pertinent word, to form a group
- 7 of citing text units,
- 8 - identification of the candidate text units containing at least one
- 9 pertinent word also cited in the citing text units, to form a group of co-cited text
- 10 units,
- 11 - assigning to the co-cited text units a relevance score as a function of
- 12 said citations.

13 The invention also proposes a method for determining relevance scores of text
14 units such as phrases in a textual document, characterized in that it comprises the
15 following steps:

- 16 - decomposition of the document into a plurality of text units,
- 17 - selection of at least one relevant text unit and of candidate text
- 18 units,
- 19 - determination of the set of pertinent words contained in the relevant
- 20 text unit (or units) and in each of the candidate text units,
- 21 - for each pertinent word contained in the relevant text unit (or units),
- 22 identification of the candidate text units comprising this pertinent word, to form a
- 23 group of cited text units,
- 24 - identification of the candidate text units containing at least one
- 25 pertinent word also cited in the cited text units, to form a group of co-citing text
- 26 units,
- 27 - assigning to the co-citing text units a relevance score as a function
- 28 of said citations.

29 The invention also proposes a method for determining scores allocated to words or
30 groups of words contained in text units such as phrases in a textual document,
31 characterized in that it comprises a step which consists in adding up the relevance
32 scores, determined by one of the methods above, of the text units in which said
33 words are located.

34

35 Brief description of the drawings

36 Figures 1 to 7 of the appended drawings illustrate various steps implemented in the
37 present invention.

1

2 Detailed description of preferred embodiments3 *Lexicon*

4 Resource (or element): Information resource such as a web page, a part of a web
5 page, a document, or an XML element. Each resource may itself consist of
6 resources, thus forming a tree structure.

7 Current resource: Resource accessed by the user at the current moment during
8 browsing (it is in particular the web page displayed in the main window of the
9 browser).

10 URI (Uniform Resource Identifier): Resource address. Will sometimes be used as
11 a synonym for URL (universal resource locator).

12 Link: URI placed in a resource. In general, by clicking on a link, the user can
13 access the resource pointed at by it.

14 Cite (a first resource cites a second resource): the first resource possesses a link to
15 the second resource.

16 Popular: Said of a resource which is accessed by a large number of users (for
17 example on the web) from its URI.

18 Private resource: Resource that is not accessible by a large number of users (in
19 particular which is not published on the web or is not widely known).

20 Associative storage: Addition of a link to a first resource, on a second resource, so
21 as to be able to retrieve the first resource via the associative search method.

22 Associative search: In order to retrieve a first resource, access to a relevant
23 resource with respect to a second resource to which a link to the first resource has
24 been added.

25 Added link: URI inserted by the user into a set of associated links.

26 Proposed spot: Spot presented by the system by priority since it comprises the
27 associated links that are most relevant with respect to the current context.

28 Spot: A spot is composed:

29 - of a set of links, in general associated with a reference resource. The
30 resources pointed at by the associated links are accessible (for example on the
31 web) from their respective URIs. The associated links are composed of given
32 associated links and of completed associated links,

33 - and (optionally) of one or more set of link(s) (in particular links added by
34 the creator of the spot and links added by users of the spot), proposed to the user

- 1 within the framework of the associative search method,
2 - and (optionally) of a link to said reference resource, said associated links
3 being selected as being relevant with respect to this reference resource,
4 Domain of relevance of a spot: set of resources designated by associated links of
5 this spot.
6 Given associated links: Associated links specified explicitly (by whoever creates
7 or publishes the resource with which said set is associated, or else by whoever
8 creates a spot for this resource).
9 Completed associated links: Associated links determined automatically (in
10 particular by means of a relative distillation algorithm described in the present
11 description).
12 Associated link score: Score of relevance with respect to the set of given
13 associated links. This score may be calculated by a relative distillation algorithm
14 such as one of those described in the present description.
15 Authority score: Relevance score of a resource with respect to a set of given
16 associated links.
17 Hub score: Relevance score of a resource citing other resources, representing the
18 relevance of the cited resources with respect to a set of given associated links.
19 Non-contextual score: Context independent relevance score.
20 Contextual score: Context dependent relevance score.
21 Noncontextual spot: With respect to a resource (or to a set of resources) in
22 question: Spot whose associated links comprise the URI of the resource in question
23 (or at least some of the URIs of the resources in question) with a score (or a mean
24 score) that is greater than a given threshold or that is selected in such a way as to
25 maximize it (cf. the spot selection procedure described in the present description).
26 Contextual spot: Spot whose associated links are the most relevant with respect to
27 the context.
28 Context: Browsing context.
29 Spot server: Server on the Internet providing the association between associated
30 link and spot.
31 Current spot server: Spot server to which the user is directly connected.
32 Relevant region of a resource: Part of a resource containing at least one relevant
33 link and containing no nonrelevant link.

1 *Methods of associative storage and associative search*

2 [Vocabulary used:

3 First page = page stored by the user so that he can retrieve it easily;

4 Second page = page used by the user as storage medium (to store an association
5 with the first page, which we shall subsequently refer to as "for storing the first
6 page" for the sake of conciseness);

7 Current page = page presented at the current moment in the main window of the
8 Internet browser.

9 These are for example web pages, however the first page may be a private resource
10 such as a document (text, multimedia or other document) which belongs to him].

11 The system allows the user to add a link to a first page on any second page
12 whatsoever (or in the vicinity of the second page; we shall subsequently use the
13 expression "on the second page" for the sake of conciseness).¹

14 The user accesses the pages by means of a browser furnished with the system
15 specific extension (or via an intermediate web server). Adding a link can be done
16 for example by a drag and drop: the user grabs a handle representing the first page
17 and drops it onto the second page; for example the link added is then presented by
18 the system as a vignette in the style of a "post-it" in the place where it was
19 dropped, or in a window adjacent to the main window of the browser (or in a frame
20 adjacent to the frame presenting the original web page). He can also drop it on an
21 icon representing the second page (for example in his favorite links). The system
22 then stores the relation with the user considered, the association between the link
23 on the first page and the second page in question.

24 Thereafter, when the user accesses a page relevant with respect to the second page
25 (or the second page itself), the URI² of this added link to the first page is
26 automatically presented to him.

27 Thus, to retrieve the first page, the user merely has to access any page whatsoever³
28 which is relevant with respect to the second page.

¹ The step consisting in adding a link in this manner, on a second resource, to a first resource (so as to be able to retrieve it by the method described in this report) is called associative storage.

² As well as optionally other indications pertaining to the link added, such as the text or the graphical object which accompanies the added link, or else a simplified or miniaturized presentation of the first page itself.

³ Said any page whatsoever is already or will have to be taken into account by the system. The user will thus prefer to choose a popular page to speed up the search. The system is furnished with a crawler the aim of which is precisely to take into account as many accessible pages (especially on the Internet) as possible which are of interest to the user.

1 More simply, in so far as:

2 - the user chooses said second page because it is relevant with respect to the
3 first page

4 - and that the relevance relation is transitive at this level,

5 to retrieve the first page, the user merely has to access any page (accessible by the
6 system) which is relevant with respect to the first page: this is the associative
7 search method.⁴

8 Note that during the step of associative storage the user can increase his chances
9 by adding a link to the first page on several second pages.

10 Furthermore, in so far as the relevance relations are symmetric, the added links are
11 implicitly bi-directional. Furthermore, in the case where the current page is a
12 private resource, the system can liken it to the second page(s) on which, as
13 appropriate, the user had added a link to this private resource, and present the other
14 first pages that he also added on this (these) second page(s).

15 The step of associative storage can be automated (or be computer aided).
16 Specifically, the addition of a link to a first page on a second page can be (semi-)
17 automated according to the following steps:

18 I - determine key words or main phrases of the first page (that are contained in the
19 page or associated with it - for example are delimited by "meta-tags"),

20 II - provide these key words or main phrases to a search engine which will return a
21 set of links on pages containing these key words,

22 III - take at least one subset thereof (for example the best N according to the search
23 engine) so as to use them as second pages,

24 IV - add a link to the first page on these second pages.

25 Note that as regards step I, various techniques for automatically extracting key
26 words or main phrases of a text already exist.

27 The key words may also be extracted from the text in the following manner:

28 - for each word, determine the score of this word by adding up the scores of
29 the phrases in which it is located and then normalizing these scores (for example
30 by dividing each score thus obtained by the square root of the sum of the squares
31 of all the scores thus obtained);

⁴ To facilitate reading, the storage/associative search method is described here while speaking of pages, but the method applies more widely to resources (not only to pages).

1 - select the words having the largest scores as key words.

2 The two methods presented above may be combined by retaining from the key
3 words selected only those which are located in the phrases selected. The complete
4 method for extracting the key words from the text is then as follows:

5 - remove the nonpertinent words from the text (called "stop words" in the
6 literature);

7 - identify the links between the phrases: a phrase is linked to another when it
8 contains at least one word in common;

9 - apply the absolute distillation procedure (described later), or an equivalent
10 procedure utilizing a graph of links (such as PageRank), to the phrases thus
11 interlinked, to determine their scores;

12 - for each word, determine the score of this word by adding up the scores of
13 the phrases in which it is located and normalized;

14 - select the phrases having the largest scores as being the main phrases of the
15 text.

16 As a variant, in so far as (one or) certain phrases of the text may be labeled as
17 being prevalent, to determine the scores of the phrases, instead of the absolute
18 distillation procedure it is possible to use the relative distillation procedure
19 (described later) to determine the relevance score of the phrases with respect to
20 said prevalent phrases.

21 Moreover, instead of actual phrases, it is possible to consider any kind of text parts
22 or units. The method using relative distillation thus consists in determining
23 relevance scores of co-cited "text units" (such as phrases):

24 The text units comprising at least one word in common with the prevalent unit (or
25 set of units) are identified so as to form a group of citing text units. A link is
26 created (temporarily) from each citing text unit to the prevalent text unit (or set of
27 units).

28 The text units containing at least one word also contained in the citing text units
29 are identified so as to form a group of co-cited text units. A link is created
30 (temporarily) from each citing unit to each co-cited unit with which said citing unit
31 possesses at least one word in common.

32 One of the methods, described later, of calculating relevance scores by the relative
33 distillation procedure is then applied. The whole set of identifiers of the relevant

1 text units constitutes the URIs of the query.⁵

2 The implementation of the associative search system will now be described.

3 To present, to a user who accesses a current page, links on first pages, the system
4 performs the following steps:

5 Step a: determine the relevance score of second candidate pages with respect to the
6 current page⁶,

7 Step b: select the (or a certain number of) second pages having (as appropriate) a
8 sufficient relevance score,

9 Step c: present to the user the (URIs of the) first pages of the links that he had
10 added on the second pages which have been selected in step b; optionally also
11 present the (URIs of the) second pages themselves to him.⁷

12 As a variant, during the associative storage, instead of adding on the second page a
13 link to the first page, the user can overlay onto the second page or insert therein
14 an annotation (or any resource such as an icon or other graphical object), which
15 then plays the role of first page within the sense of the present method. In this case,
16 during step c) of the associative search, the system presents the second page or
17 pages which have been selected while also presenting their annotations (or the
18 resource that has been added to them).⁸

19 To facilitate reading, the following 7 steps (see Fig. 1) will be considered:

- 20 • R consists of the pages⁹ of the query.
- 21 • R^- is the set of pages which contain a link to¹⁰ at least one page of the
22 query.
- 23 • R^{++} is the set of pages pointed at (cited) by at least one page of R^- .
- 24 • R^{+-} is the set of pages which cite at least one page of R^+ ($R^- \subset R^{+-}$):
- 25 • R^+ is the set of pages cited by at least one page of the query (R).
- 26 • R^{+-} is the set of pages which cite at least one page of R^+ .
- 27 • R^{++} is the set of pages cited by at least one page of R^{+-} ($R^+ \subset R^{++}$).

⁵ The set of identifiers of the citing text units constitutes the set R^- . The set of identifiers of the co-cited text units constitutes the set R^+ , and so on and so forth.

⁶ This step is composed of step a and/or step a' (see later...)

⁷ To do this, as already mentioned, the system possesses in memory the relation between user, second page (on which the user in question has added links) and first page (link added by the user in question on the second page in question). Thus the system can firstly determine the set of second candidate pages for the current user so as to perform step a, then in step c retrieve the added links to be presented to the user.

⁸ In the remainder of the description, the expression link added on a second page is understood to mean that we include this typical case where there is a resource added to the second page.

⁹ ("page" is understood to mean "page URI")

¹⁰ (stated otherwise "which cite", or else "which point at")

1 To determine the relevance score of the second candidate pages with respect to a
2 current page R (understand R here as current resource¹¹), the system implements a
3 method of "relative distillation" comprising at least one out of the following steps a
4 and a'.

5 Step a:

6 Step a-1: Identify the set R^- of pages which possess at least one link to R ,¹²

7 Step a-2: Retrieve in memory the set of second candidate pages for the current user
8 and perform the intersection between the set R^{++} of the pages pointed at by the
9 pages of R^- (note that R is in the set R^{++}) and the set of second candidate pages for
10 the current user;

11 Step a-3: For each page of the set resulting from step a-2, calculate its relevance
12 score (authority score) with respect to R . (Note that this step includes the
13 identification of the set of pages of R^{++} possessing at least one link pointing to at
14 least one subset of the set resulting from step a-2 - see the "selection of spots"
15 section).

16 Step a':

17 Step a'-1: Identify the set R^+ of pages pointed at by R ;

18 Step a'-2: Retrieve in memory the set of second candidate pages for the current
19 user and perform the intersection between the set R^{++} of pages possessing at least
20 one link to a page of R^+ (note that R is in the set R^{++}) and the set of second
21 candidate pages for the current user;

22 Step a'-3: For each page of the set resulting from step a'-2, calculate its relevance
23 score (hub score) with respect to R . (Note that this step includes the identification
24 of the set of pages of R^{++} pointed at by at least one subset of the set resulting from
25 step a'-2).

26 The calculation of the relevance scores in steps a-3 and a'-3 may be performed by
27 means in particular of one of the equations presented later in the "selecting the
28 spots" section which moreover describes improvements to the method presented
29 above. In particular the scores are sharpened by successive iterations. During these
30 iterations, the hub pages in step a and the authority pages in step a' also acquire
31 relevance scores (hub scores and authority scores respectively). In addition to the
32 second candidate pages (that is to say in addition to the URIs of the pages of R^{++} in
33 step a and/or of R^{++} in step a') determined as described hereinabove, it is then also

¹¹ Since here the query is formed of a single page.

¹² A web search engine can be used to determine the resources that point to a given resource.

1 possible to include, in the resulting set provided at step b, the hub pages of step a
2 and the authority pages of step a' (since they now have relevance scores).
3 Moreover the weights of the links between close pages¹³ are diminished so as to
4 further improve the results.

5 The system can therefore select the second pages that are most (or sufficiently)
6 relevant to step b and perform step c to present their added links to the user.

7 The results obtained by the relative distillation method may be stored (then
8 maintained - see later the "maintaining the spots" section) with the aim of avoiding
9 recalculating them during accesses to the current pages already processed. Thus,
10 the system maintains, in a second memory, the scores of the second pages with
11 respect to the current pages in the cases where these scores are greater than a given
12 threshold. For a current page already processed, the response of the system is then
13 almost immediate.

14 Stated otherwise, step a is modified as follows:

15 Step a': Consult the second memory to ascertain whether the second pages most
16 relevant for the current page have already been stored (and if these data in memory
17 are sufficiently fresh), as appropriate go to step c, otherwise determine and store
18 the relevance score of second candidate pages with respect to the current page.

19 As a variant, the system stores (then maintains - see later the "maintaining the
20 spots" section) the necessary data without waiting for a user to access a current
21 page; storage is triggered by the use, by the user, of a new second page (as
22 associative storage medium).

23 By utilizing the fact that the relevance scores are reflexive¹⁴, the system starts from
24 each second page to construct R^- and R^{++} and (R^{-+}) and/or R^+ and R^{+-} (and R^{++}),
25 calculates by relative distillation the relevance scores of all the potential current
26 pages, and stores them in a second memory (this being an inverse memory able to
27 provide, for each potential current page, the second relevant pages).

28 Moreover, as already indicated, the system maintains a first memory containing the
29 links added by user and second page.¹⁵

30 Thus, when a user actually accesses a current page, the system selects from the

¹³ To identify the closeness of the pages to the ends of the links the system additionally identifies the set of pages R^- of the pages possessing at least one link to the pages R^- and the set of pages R^{++} of the pages possessing at least one link to the pages R^{++} (see the "filtering" section).

¹⁴ (i.e. the relevance score of a second page with respect to a current page is equal to the relevance score of this current page with respect to this second page)

¹⁵ Note that, advantageously, the data in the second memory are not per user and may thus serve all the users.

1 second memory the second pages - from among the second pages used by this user
2 as storage medium¹⁶ - which have the highest relevance scores with respect to said
3 current page, then retrieves (from the first memory) the links added by this user on
4 these second pages.

5 Stated otherwise, the method comprises the following steps¹⁷.

6 For each new second page R (on which a user adds a link)¹⁸:

7 Step m1: Perform at least one of steps m1-1 and m1-1', then perform step m1-2:

8 Step m1-1:

- 9 - identify the set R^- of pages which possess at least one link to R;
- 10 - identify the set R^{++} of potential current pages pointed at by the pages of R^- ;
- 11 - for each page of R^{++} (except R) calculate its relevance score (authority
12 score - see the "selecting the spots" section) with respect to R; note that this step
13 includes the identification of the set of pages R^{+-} possessing at least one link
14 pointing at at least one subset of R^{++} (see the "selecting the spots" section);

15 Step m1-1':

- 16 - identify the set R^+ of pages to which R possesses at least one link;
- 17 - identify the set R^{+-} of potential current pages pointing to at least one page
18 of R^+ ;
- 19 - for each page of R^{+-} (except R) calculate its relevance score (hub
20 score - see the "selecting the spots" section) with respect to R; note that this step
21 includes the identification of the set of pages R^{++} pointed at by at least one subset
22 of the elements of R^{+-} ;

23 Step m1-2: store, in a second memory, the URIs of the pages having a sufficient
24 relevance score with respect to R, in relation to R, in such a way that on the basis
25 of the URI of each of said pages having a sufficient relevance score with respect to
26 R it is possible to retrieve¹⁹ (the second page) R as well as said sufficient relevance
27 score;

28 Step m2: (in parallel with step m1) store in a first memory, for each user and each
29 second page, the added links that said user has added on said second page;

30 During access to a current page by a user:

¹⁶ (they are indicated in the first memory)

¹⁷ Steps m1 and m2 describe the associative storage method, steps a, b and c describe the associative search method.

¹⁸ Step m1 is performed only for the new second pages, while step m2 is performed each time a second page is used by a user, whether or not it is new for the system.

¹⁹ (As well as the other second pages, as appropriate, for which the relevance score of R is sufficient)

1 (Step a is no longer necessary since the scores are already in memory).

2 Step b-m: Select from the second memory a certain number of second pages²⁰,
3 from among the second pages used by said user (that are indicated in the first
4 memory), for which the relevance scores of said current page are the highest (if
5 they exist);

6 Step c (unchanged): retrieve from the first memory the links added by said user on
7 the second pages selected in step b-m and present them to said user (with
8 optionally the second pages on which they have been added and in a sorted
9 manner).

10 The improvements presented later in the "selecting the spots" section will also be
11 applied. In particular as the scores are sharpened by successive iterations, the hub
12 pages in step m1-1 and the authority pages in step m1-1' also acquire relevance
13 scores (hub scores and authority scores respectively) and may thus be included in
14 the resulting set provided in step m1-2 (in addition to the URIs of the pages of R^{+}
15 in step m1-1 and/or of R^{+} in step m1-1'). Moreover, here also the weights of the
16 links between close pages are diminished so as to improve the results (see the
17 "filtering" section).

18 With this latter method, the added links are presented almost immediately by the
19 system in all cases, that is to say even when a current page is accessed by a user for
20 the first time.

21 It was mentioned that during the associative storage step the user can increase his
22 chances by adding a link to the first page on several second pages. He will now be
23 allowed to form groups of second pages to which is added a link to the first page
24 (the idea being that, as the first page may be of interest with respect to more than
25 one center of interest of the user, the groups make it possible to class the first page
26 with respect to distinct centers of interest, each group corresponding to a different
27 center of interest).

28 Specifically, each time the user adds a link (to the first page) on a new second
29 page, the group or groups of second pages that he had already formed, as
30 appropriate, for the first page are proposed to him by the system and he can then
31 choose one or more of these groups into which to insert said new second page, or
32 otherwise he can create a new group formed of the single new second page.

33 At the same time he can also manipulate his groups more widely, such as for
34 example delete a second page of a group, split a group into two, merge two groups,

²⁰ Normally, in the second memory, the URIs of the second relevant pages with respect to a potential current page are already sorted by relevance score.

1 delete a group, etc. Finally, he can also duplicate a group so as to add thereto a link
2 on another first page.

3 Each group is processed by the system as a relative distillation query. In a similar
4 manner to the last method described²¹, for each query R (that is to say for each
5 group of second pages) the system identifies and stores (then maintains – see later
6 the "maintaining the spots" section) the potential current pages which have a
7 sufficient relevance score, and thus forms an inverse memory able to provide, for
8 each potential current page, the most relevant queries (that is to say the most
9 relevant groups).

10 Stated otherwise, the associative storage comprises the following steps:

11 (Step m1 is performed only for the queries not already known by the system or not
12 sufficiently fresh, while step m2 is performed for all the users' queries, whether or
13 not they are new for the system).

14 Step m1: Perform at least one of the steps m1-1 and m1-1', then perform step m1-
15 2:

16 Step m1-1:

- 17 - identify the set R^- of pages which possess at least one link to a page of R ;
- 18 - identify the set R^{+-} of pages (seen as potential current pages) pointed at by
19 at least one page of R^- ;
- 20 - for each page of R^{+-} (except R) calculate its relevance score (authority
21 score - see the "selecting the spots" section) with respect to R ; note that this step
22 includes the identification of the set of pages R^{+-} possessing at least one link
23 pointing at at least one subset of R^+ (see the "selecting the spots" section);

24 Step m1-1':

- 25 - identify the set R^+ of pages to which at least one page of R possesses at
26 least one link;
- 27 - identify the set R^{+-} of potential current pages pointing to at least one page
28 of R^+ ;
- 29 - for each page of R^{+-} (except R) calculates its relevance score (hub score)
30 with respect to R ; note that this step includes the identification of the set of pages
31 R^{+-} pointed at by at least one subset of R^+ ;

32 Step m1-2: Store, in a second memory, the URIs of the pages having a sufficient
33 relevance score with respect to R , in relation to R , in such a way that on the basis
34 of the URI of each of said pages having a sufficient relevance score with respect to

²¹ The difference is that here R represents a query formed of one or more resources whereas before R represented a single resource (a single second page).

1 R it is possible to retrieve²² R as well as said sufficient relevance score;

2 Step m2: (in parallel with step m1) store in a first memory, for each user and
3 query, the added links (to first pages);

4 During access to a current page by a user:

5 Step b-m: Select from the second memory a certain number of queries, from
6 among the queries (groups) used by said user as associative storage medium (that
7 are indicated in the first memory), for which the relevance scores of said current
8 page are the highest (if they exist);

9 Step c: retrieve from the first memory the links added by said user on the queries
10 selected in step b-m and present them to said user, with optionally:

11 - the (or a certain number of the) queries on which they have been added,
12 - as well as a certain number of (links to) relevant pages having a relevance
13 score estimated (in step m1-2) to be sufficient with respect to said queries selected
14 in step b-m.²³

15 The improvements presented later in the "selecting the spots" section will also be
16 applied. In particular as the scores are sharpened by successive iterations, the hub
17 pages in step m1-1 and the authority pages in step m1-1' also acquire relevance
18 scores (hub scores and authority scores respectively) and may thus be included in
19 step m1-2 (in addition to the URIs of the pages of R^{++} in step m1-1 and/or of R^{+-} in
20 step m1-1'). Moreover, here also the weights of the links between close pages are
21 diminished so as to improve the results (see the "filtering" section).

22 In step b-m, the system provides a set of selected queries. It would be
23 advantageous to sharpen the selection in such a way as to present to the user (the
24 request or requests²⁴ that are the most relevant with respect to the user's browsing
25 context. This is what will now be described.

26 The history of a user's browsing is modeled with the aid of a "context stack",
27 where with each link (that may be presented to the user) is associated a relevance
28 score at each browsing level, and when a link is nonexistent it is likened to a link
29 whose score is equal to zero.

30 When the user clicks on a link and accesses a new page, the system adds a level to
31 the context stack. On the other hand, when he clicks on the "back" command of his

²² (From among the set of queries stored, as appropriate, for this page)

²³ These URIs are analogous to "related links" mentioned in the "state of the art" section, however they are more relevant since their relevance scores have been calculated with respect to the query with which they are associated by relative distillation.

²⁴ (With the first pages and the corresponding relevant links)

1 browser the system pops a level.

2 For a given link, the contextual score is an average of the noncontextual scores²⁵ at
3 each level of the context stack, these scores being weighted as a function of depth.
4 So as not to have to recalculate all the scores each time, an exponential weighting
5 is used, this implying that the contextual score at a certain level is the weighted
6 average of the noncontextual score at this level and of the contextual score at the
7 previous level.

8 Stated otherwise, for a given URI, s being the noncontextual score at the last level
9 and r the contextual score at the previous level, the contextual score at the last
10 level is: $\lambda.r + (1 - \lambda).s$ (λ being a constant weighting between 0
11 and 1, in principle less than 1/2: the larger λ is, the more important is the
12 past).

13 Among the queries (that is to say the groups) selected in step b-m, the system
14 selects those which are closest to the context, that is to say those for which the
15 scores of the URIs stored in step m-2 are the closest to the contextual scores for the
16 user in question. To determine the closeness of each request with the context, the
17 system calculates the sum of the products, for each URI of the query, of the
18 (noncontextual) score of the query with the contextual score for the user in
19 question.

20 Step b-m is thus replaced by the following step b'-m:

21 Step b'-m: select from the second memory a certain number of queries, from
22 among the queries (groups) used by said user as associative storage medium (and
23 indicated in the first memory), for which the relevance scores of said current page
24 are the highest (if they exist) and for which the relevance scores of the potential
25 current pages are the closest to the contextual relevance scores.

26 We shall now describe a method, utilizing the system of cookies, for recognizing
27 the user when he goes from one site to another, in such a way as to be able to
28 maintain his context stack.

29 Let us recall that the cookies system allows servers of sites of an Internet domain
30 (i.e. domain name or IP address) to recognize a user (that is to say his computer)
31 when he accesses web pages belonging to one and the same Internet domain.

32 The method described here allows a server, which implements our method – it will
33 be called a client server (CLI) – to recognize even users who browse from one site
34 to another which do not form part of one and the same Internet domain, even

²⁵ (That is to say determined taking no account of the context)

1 though in their browsing these users pass through sites that do not implement our
2 method.

3 To do this, three communication mechanisms are used:

4 1 – Each web page of a site of a client server contains a frame whose address is
5 that of a centralized server (URS) which manages our method of recognizing the
6 user (USER);

7 2 – The centralized server and each client server each have a cookie stored in the
8 user's computer (note that the creation time for these cookies may be used to
9 estimate the reliability of recognition of the user);

10 3 – The client server communicates with the centralized server directly.

11 There are three possible cases which are described hereinafter (see figure 2).

12 New user for the client server and for the centralized server:

13 1. The user (the USER computer) opens a page of the clients site (CLI
14 server); there is no CLI cookie.

15 2. CLI asks URS for a free identifier for USER and receives ID="123456"

16 3. CLI sends back a page comprising two frames to USER

17 • the first frame is at the address [http://URS.com/...?ID="123456"](http://URS.com/...?ID=)

18 • the second frame is at the address <http://CLI.com/...>

19 4. USER sends the http query to URS to ask for the content of the first frame
20 ([http://URS.com/...?ID="123456"](http://URS.com/...?ID=)); as there is no cookie belonging to URS, URS
21 concludes that this is a new user and allocates him the identifier "123456".

22 5. URS responds and installs a cookie (containing ID="123456" at USER

23 6. (In parallel with 5.) URS transmits [ID="123456" (no replacement)] to CLI

24 7. (In parallel with 4.) USER sends CLI the http query to ask for the content
25 of the second frame

26 8. (After receipt of the identifier at point 6) CLI sends USER the content of
27 the frame <http://CLI.com/...>

28 New user for the client server but not for the centralized server:

29 1. USER opens a page of the client site (CLI server); there is no CLI cookie.

30 2. CLI asks URS for a free identifier for USER and receives ID="123456"

31 3. CLI sends back a page comprising two frames to USER

32 • the first frame is at the address [http://URS.com/...?ID="123456"](http://URS.com/...?ID=)

- 1 • the second frame is at the address `http://CLI.com/...`
- 2 4. USER sends the http query to URS to ask for the content of the first frame
- 3 (`http://URS.com/...?ID="123456"`) as well as the content of the cookie (created
- 4 during a previous access and comprising the identifier `ID="ABCDEF"`)
- 5 5. URS responds
- 6 6. (In parallel with 5.) URS transmits [`ID="ABCDEF"` replacing
- 7 `ID="123456"`] to CLI (+ optionally extra data specific to `ID="ABCDEF"`)
- 8 7. (In parallel with 4.) USER sends CLI the http query to ask for the content
- 9 of the second frame
- 10 8. (After receipt of the identifier "ABCDEF" at point 6.) CLI sends USER the
- 11 content of the frame `http://CLI.com/...` as well as a new cookie comprising
- 12 `ID="ABCDEF"` as replacement for the previous one
- 13 User already known to the centralized server and to the client server:
- 14 1. USER opens a page of the client site (CLI server) and transmits the content
- 15 of the cookie associated with CLI (`ID="ABCDEF"`)
- 16 2. (This step is not applicable)
- 17 3. CLI sends back a page comprising two frames to USER
- 18 • the first frame is at the address `http://URS.com/...?ID="ABCDEF"`
- 19 • the second frame is at the address `http://CLI.com/...`
- 20 4. USER sends URS the http query (`http://URS.com/...?ID="ABCDEF"`, to
- 21 ask for the content of the first frame) as well as the content of the cookie (created
- 22 during a previous access and also comprising `ID="ABCDEF"`)
- 23 5. URS responds
- 24 6. (Optionally, CLI can ask for and/or receive extra data from URS for
- 25 `ID="ABCDEF"`)
- 26 7. (In parallel with 4.) USER sends CLI the http query to ask for the content
- 27 of the second frame
- 28 8. CLI sends USER the content of the frame `http://CLI.com/...` (as appropriate
- 29 after receipt of the data in step 6.)
- 30 The method described above makes it possible to select the links to be displayed in
- 31 the web pages as a function of the browsing context²⁶. This is what will now be
- 32 described.
- 33 Let us start from the situation where each query (the server which hosts it)
- 34 possesses a set of initial URIs as well as the set of links that could be proposed to

²⁶ (Or, as described above, to select the queries themselves; this being trivial, it is not described again)

1 the user with their default scores: the noncontextual scores.

2 As already described, the contextual score is an average of the noncontextual
3 scores, weighted as a function of depth, at each level of the context stack. Thus, r_i
4 being the noncontextual score at the last level and \tilde{r}_i the contextual score at the
5 previous level, its value after having followed a link is: $\tilde{r}_i \leftarrow \lambda \tilde{r}_i + \bar{\lambda} r_i$ ²⁷.

6 The links presented to the user are those which have the largest contextual score.

7 The context stack can be displayed in the URS frame (the first frame) introduced
8 above. Thus the user can see which pages are the ones that were involved in the
9 calculation of the pages to be displayed. He can click elements of the stack to
10 climb back up the levels, and an "Erase" button makes it possible to empty the
11 context stack.

12 The context stack is stored, for each user, in the centralized server (URS), with the
13 user's identifier. Thus, each time a user opens a page at a client server (CLI), the
14 latter, having obtained the user's identifier, will give URS the noncontextual
15 scores²⁸, which will respond with the contextual scores after having performed the
16 weighted average described above²⁹. The server of the client site may then display
17 in the page the links which have the best score.

18 The steps are thus as follows (see figure 3):

- 19 1. The user (USER) sends an http query to open a page.
- 20 2. The client server (CLI) transmits the noncontextual scores of the page in
21 question and the user's identifier to the centralized server (URS)
- 22 3. URS adds a level to the context and calculates the contextual scores
- 23 4. The contextual scores (at least the best of them) are returned to the client
24 server
- 25 5. The client server selects the links which have the best score and presents
26 them to the user.

27 It may be beneficial on the one hand to group the links in various parts of the
28 pages, or even to hierarchize the parts, that is to say to allow parts to contain

²⁷ Thus giving $\tilde{r}_i = \bar{\lambda} \sum_{n=0}^{d-1} \lambda^n r_{i,n} + \lambda^d r_{i,d}$ with d the depth of the root and $r_{i,n}$ the score of page P_i

at depth n .

²⁸ To avoid unnecessary traffic it is possible to select the pages to be sent, taking only those that have a score greater than a certain threshold, for example half the threshold required in order for a page to be displayed to the user

²⁹ This is performed within the framework of step 6 described above.

subparts, in addition to links. Here are the changes that this involves:

- The current context³⁰ must contain context information for each part of the page displayed, hence when the page sends its noncontextual scores, it sends as many of them as there are parts, and URS responds to it with a context for each part. To avoid certain problems (see the following points), a default context is also necessary, representing the page itself and its parts and aggregating all the scores of all the links

- When the user clicks on a link, the context of the part which contains this link must be used as last-level context (i.e. that context will be used for the calculation of the scores at the subsequent levels). A means of obtaining this result is to place in the addresses of the links an argument which contains an identifier (unique for the page) of the part, which identifier is also transmitted to URS with the noncontextual scores.

- In the implementation of the method described here, care must be taken not to confuse the parts of various pages, for example if the user has opened several windows of his browser and clicks in a window after having clicked in another (URS stores only a context stack). This may be done by comparing the field HTTP Referer with the address of the last level of the stack and take no account of the part number other than in the case of equality. In other cases (also if the user has passed through a page of a nonclient site), the default context is taken.

A more complete example (see figures 4 and 5):

Here therefore is what happens when the user, already in a particular context (for the page `cl/com/main.html`), clicks on a link <http://CLI.com/index.html?part=1> (part = 1 signifies that the user has clicked in part 1). It is assumed that the client server CLI does not yet know the user:

(1) The browser (USER) sends the query <http://CLI.com/index.html?part=1> to the server of the client site (CLI), additionally giving him the Referer <http://cl.com/main.html> (the address of this frame).

(2) CLI will ask URS for a free number (it responds to it with 12345) for this user

(3) CLI responds to (1) with a page comprising two frames whose addresses are <http://URS.com/default.html?newID=12345> and <http://CLI.com/main.html> respectively. He also gives him a temporary cookie (session cookie) newID=12345.

(4) The user being known to URS, it has a cookie with its true identifier (678910). By loading the frames, it (its browser) will send a query for the page

³⁰ That is to say the set of contextual scores of the links at the current level.

1 <http://URS.com/default.html?newID=12345> with the cookie ID=678910.

2 (5) The user also sends a query for the page <http://CLI.com/main.html> with the
3 session cookie newID=12345.

4 (6) Having received (5), the client CLI sends URS its address
5 (<http://CLI.com/main.html>), its noncontextual scores, for each part of the new
6 page, the identifier newID=12345, as well as the part number (part=1) that it
7 received to the message (1).

8 (7) When it has received (4) and (6), URS looks at the context of the user for
9 part 1, verifies that the source page (<http://CLI.com/main.html>) corresponds to the
10 last level of the context stack for this user (otherwise it would have ignored the
11 part number and taken the default part ("D"). Thereafter it calculates, for each part
12 of the new page the new contextual scores.

13 (8) URS, having received the message (6), can respond to the message (4) of
14 the user (presenting him with the new context stack and the <ERASE> button).

15 (9) URS also responds to the message (6) from CLI, sending it the true
16 identifier of the user (678910), as well as the contextual scores.

17 (10) CLI can now respond to the message (1), giving the user are true identifier
18 (permanent cookie ID=678910, for the site CLI.com), as well as the personalized
19 page.

20 The concept of user can in reality encompass several users who share added links
21 (and the groups which serve them as support). Of course, a finer organization of
22 the users according to the added links that they share is possible.

23 We shall now describe the case where an end user subscribes to a provider user so
24 that, according to the context, the system proposes the groups and first pages (in
25 the sense of the groups and first pages described hitherto) created by the provider
26 user to the end user. The first pages may in particular be advertisements which (by
27 virtue of the capabilities of the system as hitherto) are automatically selected with
28 respect to the context.

29 The groups created by the provider user and proposed by the system to the end
30 user are called "spot".

31 The provider user manipulates and utilizes the spots as described hitherto for the
32 groups of second pages.

33 The end user can use a spot as storage medium by making a personal version
34 thereof and adding thereto a link to a first page (this is described later).

1 The main advantage of this approach is to afford the possibility of creating new
2 spots (and the expensive calculations of scores that they involve) to certain users
3 only (namely the provider users) and to offer the function of storage/associative
4 search by way of pre-existing spots (which is not expensive in terms of machine
5 resources) to all users.

6 *Spot*

7 The system that we shall now describe provides relevant links (also known as
8 "related links", see above the "state of the art" section). However, rather than
9 searching for relevant links directly, our system searches firstly to see whether
10 there exists a spot – or reference resource – whose associated links are sufficiently
11 close to the current resource or to the browsing context of the user. If such is the
12 case, the system returns the spot(s) whose associated links are the closest, as well
13 as its associated links offered in the guise of relevant links.

14 Typically the spot is proposed in a window adjacent to the main window of the
15 browser, like the existing systems providing "related links", however in contra-
16 distinction to these existing systems

17 - the system of the invention presents relevant links determined according to
18 a relative distillation method (detailed later),

19 - the browsing context taken into account by our system is not necessarily
20 solely the current page, but may include the set of resources accessed recently by
21 the user (using the system) and which are relevant with respect to the current
22 resource³¹

23 - the spots serve as associated memory for the provider users; specifically,
24 when a spot is presented to an end user, the links to first pages (or other added
25 resources³², as described previously) added by the provider user who created the
26 spot are presented to said end user³³,

27 - the spots serve as associative memory for the end users; specifically, when
28 the end user adds a link to a first page on a second page (as described hitherto), in
29 reality he adds a link on his personal version of the spot proposed for this second
30 page or for the current context.

31 Furthermore, presenting the end user with relevant links by way of spots offers

³¹ See above the description of the method of selecting groups of second pages (here of spots) according to the user's browsing context.

³² The latter include in particular advertisements billed to promoters. Advantageously, these advertisements are relevant with respect to the context (in any event the spots which serve them as support are).

³³ (The latter possibly moreover being said provider user who created the spot)

1 advantages per se, such as prompting to click in order to access the reference
2 resource (that is to say the page presenting the spot).

3 Let us now examine a few typical storage/associative search scenarios
4 implementing spots.

5 First scenario of use:

6 The provider user creates a new resource or chooses an existing resource (for
7 example a web page which he wishes to access, or a particular element contained
8 in a page ...) so as to make thereof the reference resource of a new spot.

9 To do this, he allocates it at least one given associated link pointing to a popular
10 page.

11 The system completes the set of associated links³⁴ (as described in the "selecting
12 the spots" section).

13 Thus, in the future, each time an end user accesses a resource pointed out by one of
14 the links associated with this spot, this spot may³⁵ be proposed to him.

15 Also, as described in the subsequent two scenarios of use, end users may then use
16 this new spot as storage medium (in a manner analogous to the use of a second
17 page or of a group of second pages, described above).

18 The creator of this spot thus has the advantage not only of putting it to his own use
19 but also of seeing it proposed to end users. As a link on the reference resource
20 (prompting the user to click) is included in the presentation of the spot, the
21 reference resource is thus promoted to the end users. Moreover, its added links
22 (such as advertisements) on this spot will be presented to the end users.

23 Second scenario of use:

24 On the web the end user "lands" on a first page (or other type of resource) that is
25 so interesting that he would like to store it in order to be able to retrieve it easily
26 and land back on it spontaneously when he accesses resources that are relevant
27 with respect to it.

28 Let us assume that no spot is spontaneously proposed by the system for this page.³⁶

³⁴ This is the equivalent of the second memory described in the previous section.

³⁵ It will not necessarily be this spot that is proposed but rather, among all the spots whose associated links point to resources forming the current context, the spot in which these associated links have the highest relevance scores (or the spots in which these associated links have the highest relevance scores). The selection of the spot (or spots) is described in the "selecting a spot" section.

³⁶ In the converse case, on (his personal version of) this spot, the user will directly add a link to this first web page. Note however that this action is not strictly necessary. Specifically, already without doing anything the user will have to retrieve this first page by visiting a close page that is not very popular (in the guise of relevant link associated with this same spot or with a neighboring spot).

1 The user visits a (at least one) second page, which is relevant with respect to the
2 first,

3 - and for which he knows that a spot is proposed,

4 - or else he chooses a web page which is popular since it is thus more
5 probable that a spot is proposed for it,

6 and on the spot which is proposed for this second page he adds a link to this first
7 page (for example by selecting a graphical object representing a first page and by
8 performing a drag and drop thereon on the second page, as described at the start of
9 the description).

10 In the future, this added link will then be presented to him spontaneously each time
11 that this same spot, or that a close spot, is proposed to him for the current context
12 of his browsing.

13 Third scenario of use:

14 The end user wishes to store a private resource (such as a document which belongs
15 to him and which is not published on the web). The private resource here plays the
16 role of first page.

17 He accesses a (second) page which is relevant with respect to his private resource
18 (and which preferably is popular, or for which he knows that a spot is proposed)
19 and he adds thereto a link to his private resource (that is to say he inserts this link
20 into his personal version of the spot proposed for this second page).

21 Optionally, to reinforce his action, he will also add a link (to his private resource)
22 on yet (other spots which are proposed to him for) other second pages that he finds
23 relevant with respect to his private resource.

24 In the future, a link to his private resource will be presented to him spontaneously
25 each time that one of the spots that was proposed to him for the second page or
26 pages, or that a close spot, is proposed to him for the current context of his
27 browsing.

28 Thus, in the last two scenarios above, a link to the first page is presented to the
29 user spontaneously each time that he visits pages in the domain of relevance
30 covered by the spots proposed for the second pages³⁷.

However, by doing this action the user has the extra advantage of being able to retrieve it in the
guise of link added explicitly by him, that is to say in such a way that it is made evident.

³⁷ And insofar as the second pages were chosen by the user because according to him they are
relevant with respect to the first page, and the relevance relation is transitive at this level, a link to
the first page is presented to the user spontaneously each time he visits pages which according to
him are in the domain of relevance of the first page!

1 *Selecting the spots*

2 Before the spot(s) selection step proper, the system must obtain the set of
3 "completed associated links" from the set of "given associated links" (which are
4 given by the provider user, as described in the first scenario of use).

5 Completing the associated links:

6 The set of resources pointed at by the given associated links is the query R.

7 The calculation of the completed associated links is performed by means of the
8 "relative distillation" method, comprising the following steps:

9 Step 1: Identify the set R^- of resources which possess at least one link pointing at
10 an element of R.

11 Step 2: Identify the set R^{++} of resources pointed at by the elements of R^- (note that
12 R^{++} includes R).

13 Step 3: For each resource of R^{++} calculate its authority score with respect to R.
14 (This step can include the identification of a part of the resources of R^{++}
15 possessing a link pointing to a resource of R^-)³⁸.

16 Final step: Select the elements of R^{++} having the largest authority scores.

17 The calculation of the scores in step 3 may be performed by calculating, for each
18 resource of R^{++} , the ratio between

19 - the cardinality of the set of resources which point to it AND to the
20 resources of the query and

21 - the cardinality of the set of resources which point to it OR to the resources
22 of the query

23 (or by means of one of the more complete equations described later, see in
24 particular the equation for the quantity of common reasons – or homogeneity of a
25 set of resources).

26 The authority scores are normalized (in such a manner that their sum becomes
27 equal to 1).

28 The authority scores having been obtained, they can be put to use to allocate hub
29 scores to the elements of R^- :

30 Step 4: The hub score of each element of R^- is obtained by taking the sum of the
31 authority scores (calculated in step 3) of the elements of R^{++} to which it points. The

³⁸ The resources of R^{++} will start to be taken into account right from the first iteration, as described later.

1 hub scores are normalized (in such a way that their sum becomes equal to 1).

2 Iteration restarting from step 3: the hub scores having been obtained, they can be
3 put to use to sharpen the calculation of the authority scores. Step 3 then takes
4 account of the hub scores so as not to consider all the elements of R^- on an equal
5 footing (the resources of R^- pointing to resources having a higher authority score
6 will thus have a greater influence). The cardinalities used to calculate the authority
7 scores are thus replaced by weighted cardinalities. That is to say each hub
8 resource, instead of counting for one, counts proportionately to its hub score. (The
9 equations are detailed later).

10 Step 3 then includes the taking into account of the resources of R^{+-} pointing to the
11 resources of R^+ having the largest authority scores, in addition to R^- (a method
12 optimizing the way in which R^{+-} is taken into account is described later).

13 After step 3, we can optionally perform step 4 again, and so on and so forth until
14 convergence, that is to say until the difference between the results obtained in the
15 last iteration and those obtained in the previous iteration are negligible (in general,
16 fewer than 10 iterations are sufficient).

17 Variant for step 2: to form R^+ , instead of taking all the links contained in the
18 resources R^- the system will take only the links located in the relevant regions of
19 the resources of R^- . As these relevant regions can be determined only onward of
20 the moment at which the hub scores of the links that they contain are known, this
21 variant will be implemented only onward of the first iteration, that is to say after
22 having performed step 4 the system will iterate restarting from step 2 rather than
23 from step 3.

24 Variant for step 3:

25 With each link possessed by a resource of R^- (or of R^{+-}) is associated a weight
26 equal to the complement of the closeness of the two resources connected by this
27 link. Thus, the links connecting two close resources will be weakened. Thus the
28 importance of the links between the resources which mutually promote one another
29 (for example because it form part of one and the same web site and mutually cite
30 one another) is thus decreased. Once the links are thus weighted, the system
31 calculates the authority scores, not now by using the sum of the hub scores, but the
32 sum of the hub scores multiplied by their weights (this is detailed and illustrated by
33 an example later).

34 The closeness of the two resources connected by the link in question is obtained by
35 calculating the ratio between

36 - the cardinality of the set of resources which point to the two connected

1 resources and
2 - the cardinality of the set of resources which point to at least one of the
3 connected resources.

4 (or by means in particular of one of the more complete equations described later).

5 It is also advantageous to perform the same algorithm downstream, that is to say
6 by calculating the hub scores of the resources of R^{+-} (which downstream cite the
7 same resources as the query).

8 The downstream algorithms are identical to those upstream except that B
9 (backward) is replaced by F (forward) and vice versa³⁹, and $-$ is interchanged with
10 $+$ (e.g. R^{+-} is replaced by R^{++}).

11 Consideration will also be given, advantageously, to the hub resources upstream
12 and the authority resources downstream, in such a way that the hub pages in step
13 m1-1 and the authority pages in step m1-1' also acquire relevance scores (hub
14 scores and authority scores respectively) and may thus be included in the resulting
15 set provided at step m1-2 (in addition to the URIs of the pages of R^{+-} and/or of R^{++}
16).

17 By completing the associated links of each new query (spot) introduced, the
18 system forms an inverse memory able to provide, for each potential current
19 resource corresponding to an associated link, the most relevant queries (that is to
20 say the most relevant spots).

21 Stated otherwise, the associative storage now comprises the following steps:

22 (Step m0 is performed independently of the other steps. Step m1 is performed
23 only for the queries, not already known by the system or not sufficiently fresh,
24 introduced by a provider user, while step m2 is performed for each use of a query
25 (that is to say of a spot) as associative storage medium by a provider user or an end
26 user.)

27 Step m0: store (in a third memory) the usage rights for spots for each user.

28 Step m1:

29 Step m1-1 corresponds to completing the associated links as described
30 hereinabove.

31 Step m1-2: store, in a second memory, the URIs of the resources having a
32 sufficient relevance score with respect to R , in relation to R , in such a way that on
33 the basis of the URI of each of said resources having a sufficient relevance score

³⁹ $B(R_i)$ is the set of URIs of the pages having a link to the page R_i . $F(R_i)$ is the set of URIs of the pages to which R_i has a link.

1 with respect to R it is possible to retrieve⁴⁰ R as well as said sufficient relevance
2 score;

3 Step m2: (in parallel with step m1) store in a first memory, for each user and
4 query, the added links (to first resources);

5 During access to a current resource by a user:

6 Step b-m: Select from the second memory a certain number of queries, from
7 among the queries (spots) (indicated in the first memory) that said user has the
8 right to use, for which the relevance scores of said current resource are the highest
9 (if they exist) and for which the relevance scores of the associated links are the
10 closest to the contextual relevance scores for said user;

11 Step c: Retrieve from the first memory the links added by said user on the queries
12 selected in step b-m, as well as the links added by their creators (if they are
13 different from said user), and present them to said user, with optionally:

- 14 - the (or a certain number of the) queries on which they have been added,
- 15 - as well as a certain number of (associated links to) resources having a
16 sufficient estimated (in step m1-2) relevance score with respect to said queries
17 selected in step b-m.

18 The relative distillation method will now be detailed.

19 The essential idea of the calculation of the relevance score (of a web page P_2 with
20 respect to a given web page P_1) is as follows⁴¹:

21 Let p_1 be the probability⁴² that a random author (of a web page) places a link on P_1
22 in a page.

23 Let p_2 be the probability that a random author places a link on P_2 in a page.

⁴⁰ (From among the set of queries stored, as appropriate, for this resource)

⁴¹ Hereafter, we shall assume that P_1 and P_2 , (or P_i , P_j , etc) are web pages, although the methods described are far more general, as has already been mentioned. For example, it should be noted that instead of utilizing the hypertext links and the queries as mentioned hereinabove, the system may be based on analysis of the traces of the cutting and pasting of information fragments performed by the users (within the framework of creating and manipulating information resources), so as to automatically suggest other fragments which might enrich these resources. These traces may in fact be likened to links. For example, when part of a web page is copied into a document, the system is capable of deducing therefrom and of storing the existence in the document of a link to the web page, and the same mechanisms described here may then be applied. Moreover, the method described here may advantageously be applied by likening the links from one resource to another resource, to links from a user to a resource that he likes (that is to say to a resource which interests him). It is thus possible to determine the quantity of common reasons (between several resources) to be liked by users. This can in particular serve to categorize these resources.

⁴² The probability of being interested in a (or certain) page(s) is approximated by counting the number of pages which have a link on it (them) and by dividing this number by an estimate of the number of pages which could have had one.

1 Let $p_{1\&2}$ be the probability that a random author places a link on P_1 and a link on
2 P_2 in a page.

3 $B(P_i)$ is the set of URIs of the pages having a link to the page P_i .

4 $F(P_i)$ is the set of URIs of the pages to which P_i has a link.

5 The relevance of a page with respect to a set of pages may be defined by the
6 "quantity of common reasons" to be interested in all these pages.

7 Algebraic calculations make it possible to obtain equations giving the quantity of
8 common reasons between several pages. This quantity (or closeness, or else
9 homogeneity) is denoted x , subscripted with the pages concerned; the probability
10 of being linked to a certain page P_i is denoted p_i ; the probability of being linked to
11 at least one page out of P_i, P_j, \dots, P_n is denoted $p_{ij\dots n}$:

12 $\overline{x_{ij}} = \frac{\overline{p_i \cdot p_j}}{p_{\emptyset} \cdot p_{ij}}, \overline{x_{ijk}} = \frac{\overline{p_i \cdot p_j \cdot p_k \cdot p_{ijk}}}{p_{\emptyset} \cdot p_{ij} \cdot p_{ik} \cdot p_{jk}}$, and so on and so forth (all the
13 subsets of odd size
14 in the numerator, and the others in the denominator)⁴³.

15 This equation may be denoted more compactly thus:

$$16 \quad \overline{x_S} = \prod_{P \in S} \overline{p_P}^{\sigma_P} \text{ with } \sigma_P = (-1)^{|P|}.$$

17 The probabilities concerned above involve the number (the count) of pages of R^-
18 which contain a given link or a link from among a set of given URIs (to pages of
19 R^+). It would be beneficial to weight this number by the *quality of citation* (hub
20 score, described later) of each page which contains such a link.

21 It would thus be desirable for a page of R^- citing more better pages (of R^+) to be
22 regarded as being of better quality of citation, and for in return a higher weight to
23 be given to it within the framework of the calculation of the scores⁴⁴ of the pages
24 that it cites (R^+), the scores of the pages of R^- and those of the pages of R^+
25 mutually influencing one another in an iterative approach (bipartite reinforcement)
26 which converges⁴⁵.

27 The number of pages of R^+ citing each candidate page (that is to say of R^+) also
28 comes into the calculations. However, it is expensive to take them into account.

⁴³ The bars above indicate complements, and p_{\emptyset} , the probability of liking at least one page of an empty set, is a constant equal to zero; it is present in the equation for reasons of consistency.

⁴⁴ Recall that here one is dealing with relevance scores with respect to the query, in contradistinction to the state of the art which makes it possible to determine a score of quality "in the absolute".

⁴⁵ Note that the calculation of the relevance score of a page of R^+ may result in a negative value (that we will then neutralize; this is described later). Specifically, certain pages may not only be close to the query, but even be antagonistic with respect to it (the fact of being of interest thereto decreases the chances of liking the pages of the query and vice versa).

Hence, the results will be approximated by considering only those which cite the candidate pages having a good score, this score being calculated firstly by considering only R^- and subsequently by extending this set to R^{+-} gradually.

To calculate the relevance score of a candidate page, instead of taking the result of the equation for the quantity of reasons directly, it is preferable

- to take it together with the overall cardinalities replaced by the total of the hub scores of the pages in question and

- to multiply this result by the authority score of the candidate page (simply calculated on the basis of the total of the hub scores of the citing pages), so as thus to weaken the pages which are relatively less reliable (being less popular as they are).

After a first iteration, in the citing pages the system can

- label the regions containing directed links on pages of R^{+} having a good score

- and already begin to prune the links which are not situated in these regions.

As the links in question are located under nodes of a typically tree-like document structure (such as in HTML in particular), to determine a relevance region it suffices to take the (minimal) nodes which encompass all the good links and to take away from them the (maximal) subnodes which contain a bad link (score too low, or URI explicitly refused) and which contain no good link (sufficient score).

The algorithm makes it possible, having a homogeneous set (having sufficient homogeneity) of URIs associated with close pages, to obtain a list of URIs of pages which are relevant in regard to this set. The way in which this algorithm may be utilized to obtain a set of relevant pages for an inhomogeneous set will be described later.

As input, this algorithm takes

- a set K of reference URIs ("Kernel")
- a set A of candidate URIs ("Authority")
- a set H of hub candidate URIs
- a set T of URIs to be refused ("Trash")

We have: $K \subset H \subset A^-$ and $T \cap K = \emptyset$. (E being a set of

$$\text{URIs, } E^- = \bigcup_{P_i \in E} B(P_i) \text{ and } E^+ = \bigcup_{P_i \in E} F(P_i))$$

1. With each page P_i of H , associate a number h_i , initially set to $\frac{1}{|H|}$, its hub score⁴⁶.

⁴⁶ Thus, advantageously, the sum of the $|H|$ scores h_i is equal to 1.

1 2. (Re)calculate the authority scores:

2

3 a. For each page P_i of A , beginning with those of K , associate a number a_i , its
4 authority score, equal to

5

$$\sum_j l_{ji} \cdot h_j, \quad \text{where } l_{ji} = \begin{cases} 0 & \text{if there is no link between } P_j \text{ and } P_i \\ 1 & \text{if there is a link between } P_j \text{ and } P_i. \end{cases}$$

6

7 b. A possible but dangerous optimization: if, for certain pages, a_i is
8 sufficiently close to its value calculated previously (as appropriate), and if the
9 authority scores of the pages of K have not varied either, we can keep the old value
10 of r_i for this page, to save on calculations.

11

12 3. (Re)calculate the relevance scores:

13 a. For each page P_i of A calculate r_i^+ , equal to $w_{i \cup K}$

14 $r_i^+ = w_{i \cup K}$

15 and in the case where the result is negative (case of a page antagonistic to R)
16 neutralize the incoming links in such a way as to have $r_i^+ = 0$.

17

18 The upstream homogeneity w_S of a set S is defined as follows:

19
$$w_S = \prod_{P \in S} a_P^{\sigma_P}, \quad \text{where}$$

20
$$\sigma_P = \begin{cases} -1 & \text{if } P \text{ contains an even number of pages} \\ +1 & \text{otherwise} \end{cases}$$

21
$$a_P = \Delta \prod_i h_j l_{jp} \quad \text{where}$$

22 Δ is an arbitrary constant less than but close to 1 (it serves to avoid divisions by
23 zero but does not change the principle of the algorithm. If the set H is larger than
24 K then this constant may be equal to one

25

26
$$l_{jp} = \begin{cases} +1 & \text{if } \exists P_j \mid P \mid l_{ji} = +1 \\ 0 & \text{otherwise} \end{cases}$$

27

28
$$\text{with } l_{ji} = \begin{cases} 0 & \text{if there is no link between } P_j \text{ and } P_i \\ 1 & \text{if there is a link between } P_j \text{ and } P_i \end{cases}$$

29 Stated otherwise, l_{jp} is equal to 1 if there is a link

- 1 • from a page P_j (of H)
- 2 • to at least one page P_i of P
- 3 and zero otherwise.

4

5 This signifies quite simply that a_p is the total of the hub scores of the pages (of H)

6 which point at at least one page of P (P being the current subset of S which is

7 considered).

8 *For each existing link l_{ji} , it is possible to associate with it a weight as a function of*

9 *the closeness of the pages P_i and P_j and thus to improve the result – see later.*

10 Here, since $\forall P_i \in K$ we have $r_i^+ = w_K$ (the relevance is the same for all the pages

11 P_i of K), the relevance score r_i^+ has to be calculated only once for the pages of K

12 (besides, it will already be calculated during the procedure for chopping the query

13 R into subqueries (kernels) K , and will therefore already be known on entry to the

14 procedure).

15

16 b. (This point will be skipped the first time). To have their sum equal to 1, we

17 must divide each r_i^+ by the sum $\sum |r_i^+|$ of all the absolute values of the r_i^+ .

Let $\delta = \sum_i \left| r_i - \frac{r_i^+}{\sum_i |r_i^+|} \right|$ be the global variation of the

18

19 relevance score.

20

21 If $\delta < \varepsilon$ ($\varepsilon > 0$ being a margin of error) we assume convergence has occurred and

22 the method stops. Otherwise, the method continues.

c. We replace r_i by $\frac{r_i^+}{\sum_i |r_i^+|}$

23

$$r_i \mapsto \frac{r_i^+}{\sum_i |r_i^+|}$$

24

25 a friction factor τ also being able to be used:

$$r_i \mapsto \tau r_i + \bar{\tau} \frac{r_i^+}{\sum_i |r_i^+|} \quad (\tau \in [0;1], \text{ we shall preferably take a}$$

26

27 very small value e.g. 0.01 so that in cases where this is not necessary the number

28 of iterations does not change).

1 4. ⁴⁷For each page P_i of H :

2

3 a. Find all the links which point at a page having a relevance score larger than
4 a threshold epsilon to be chosen ($\epsilon > 0$).

5

6 b. Find I_i , the smallest HTML element⁴⁸ containing all of the links found in
7 point a above.

8

9 c. For each link pointing at a page of T (if T is not empty), find the largest
10 HTML element containing it (if there is one) and not containing any link found in
11 point a. above, and remove it from I_i .

12

13 d. We keep all the links remaining in I_i and we delete the others (or else we
14 neutralize them by setting their l_{ij} to zero).

15

16 5. Recalculate the hub scores:

17

18 a. For each page P_i of H , calculate $h_i^+ = \sum_j l_{ij} r_j$, the sum of the relevance scores
19 of the pages pointed at.

20 b.
$$h_i \mapsto \frac{h_i^+}{\sum |h_i^+|}$$

21

22 (The division by $\sum |h_i^+|$ is, as for the relevance score, so as to keep their sum equal
23 to 1).

24 Then return to point 2.

25 Initially, so as to process only a reduced number of pages, the relevance scores
26 may be calculated on the basis of R^- (if we took $H = R^-$). Hence, this will only be
27 an approximation. Specifically, for the scores to be correct, they have to be
28 calculated based rather on $H = R^{+-}$. However, as the construction of R^{+-} is
29 relatively expensive, we shall take only a subset: for R^{+-} we shall take only the
30 pages pointing at the pages of A which have a good score.

31 Thus⁴⁹, a subset will be added before the end of step 2.a:

⁴⁷ This point may possibly be ignored after the first time.

⁴⁸ (Or other analogous representation ...)

⁴⁹ Several procedures may be used; here we present the preferred one.

2.a.1. In the case where the score r_i^+ of the current page (P_i of A) is sufficient⁵⁰, r_i^+ is recalculated after having inserted the new pages of $B(P_i)$ into H .

$$H \mapsto B(P_i) \cup H$$

We introduce an authority score for the pages of A and the equation r_i^+ is $r = w_{i \cup K} \cdot a_i$ (rather than $r = w_{i \cup K}$). The new coefficient a_i will make it possible to weaken the pages that are not very reliable (because they are not very popular). Furthermore, the equation will be more consistent insofar as the relevance score will no longer be the same for all the pages of the query.

The procedure is now as follows:

1. This point is the same as that of the algorithm for calculating relevance scores presented above.
2. This point does not change either.
3. (Re)calculate the relevance scores:
 - a. For each page P_i of A calculate r_i^+ , equal to $w_{i \cup K} \cdot a_i$ and in the case where the result is negative (case of a page antagonistic to R) neutralize the incoming links so as to have $r_i^+ = 0$.
 - b. Resume from point 3.b of the previously presented algorithm for calculating relevance scores.

Filtering:

For each existing link l_{ji} , it is possible to associate therewith a weight dependent on the closeness of the pages P_i and P_j and to thus improve the result. This makes it possible to decrease the importance of the links between pages which mutually promote one another. Typically one thus succeeds in filtering for example the links of the "abstracts" and other "menus" which, repeatedly, are located in all the pages of a site.

The basic idea consists in weakening the links connecting two pages that we know to be close, by assigning a weight to each link, which weight will be equal to the complement of the closeness of the two connected pages (the greater the closeness, the more the link must be weakened). Once the links have thus been weighted, it is possible to calculate the homogeneity of a set of pages using the sum of their weights, rather than the number of citing pages.

⁵⁰ (That is to say greater than a chosen threshold; this threshold can be dependent on the current cardinality of H , specifically, the closer we get to R^{++} (e.g. H_{final}), the more chance the calculated score has of being correct)

1 In point 3.a of the algorithm, in the definition of the authority score we replace
 2 $h_j l_{jp}$ with $h_j l_{jp}$ where $l_{jp} = \min_{P_i \in P} 1; \max(\overline{l_{ji}}, \overline{x_{ji}})$

3 Explanations:

4 • $\overline{l_{ji}}$ is the complement of the closeness between page P_j and page P_i if
 5 there is a link from page P_j to page P_i , and zero otherwise

6 • $\max(\overline{l_{ji}}, \overline{x_{ji}})$ is the complement of the closeness between page $P_j \in H$ in
 $P_i \in P$

7 question and page $P_i \in P$ for which the link between P_j and P_i exhibits the
 8 minimum closeness

9 • $\min_{P_i \in P} 1; \max(\overline{l_{ji}}, \overline{x_{ji}})$ signifies that this value is truncated above to 1

10 • and always $l_{ji} = \begin{cases} 0 & \text{if there is no link between } P_j \text{ and } P_i \\ 1 & \text{if there is a link between } P_j \text{ and } P_i \end{cases}$

11 Stated otherwise, if there is at least one link

12 • from the page P_j (of H) in question

13 • to a page P_i of P ,

14 l_{jp} is equal to the complement of the closeness between page P_j and page P_i which
 15 is the least close to it and to which it possesses a link. l_{jp} is the sum of the

16 weights thus associated with the pages of H which point at at least one of the pages
 17 of the subset P considered.

18

19 To determine the closeness x_{ji} , we can take the equation (already described) for the
 20 quantity of common reasons:

$$\overline{x_{AB}} = \frac{\overline{P_A} \cdot \overline{P_B}}{\overline{P_A} \cdot \overline{P_{AB}}}$$

21

22

23 Figure 6 presents an example where the number of pages pointing at page A is
 24 equal to $0.9+0.2+0.4+0.8 = 2.3$

25 The number of pages pointing at page B is equal to $0.9+0.1+0.3+0.5 = 1.8$

26 The number of pages pointing at A or B ($N_{P_{AB}}$) is equal to
 27 $0.9+0.2+0.9+0.8+0.3+0.5 = 3.6$

28

29 Thus, if we assume that $|H| + h = 100$, the calculation of the closeness of A and B
 30 gives:

$$\overline{x_{AB}} = \frac{\overline{p_A} \cdot \overline{p_B}}{\overline{p_a} \cdot \overline{p_{AB}}} = \frac{0.977 \cdot 0.982}{1 \cdot 0.964}, \quad \text{this giving} \quad \tilde{x}_{AB} = \frac{x_{AB}}{p_B} \approx 0.264 = 26.4\%.$$

The filtering described above uses a weight $\overline{x_{ji}}$. Since we now have the scores⁵¹ of the citing pages, we can optionally improve the method by taking $\overline{x_{ji}} \cdot \overline{h_j}$ as weight (instead of $\overline{x_{ji}}$), where h_j is the score of the citing page (weakening a link) (from a citing page P_j to a cited page P_i) further when the score of the citing page P_j is low.

It should be noted that in order to calculate the closeness x_{ji} between two *connected* pages P_i and P_j , instead of using the equation for the quantity of reasons as illustrated hereinabove, it is possible to calculate the ratio between:

- the cardinality of the set of pages which point to P_i AND P_j
- and the cardinality of the set of pages which point to P_i OR P_j .

Determination of the homogeneous subsets of a query:

We provide the system with a set R of pages and possibly a set of pages R_x of pages that we do not explicitly want ($R \cap R_x = \emptyset$). The system will identify within R at least one group of "homogeneous" pages and will launch a separate sub-query on this or each group. These groups are called "kernel". To form the response, we shall then take a combination of the scores obtained. This method thus comprises the following steps:

1. For each page P_i of R , find $B(P_i)$, the set of pages citing P_i .
2. Find $R^* = \bigcup_{P_i \in R} B(P_i)$, the set of pages citing at least one page of R .
3. In the pages of R which are not yet in a kernel (at the start none is), find the page P_B having the largest set $B(P_B)$ of incoming links⁵² and create a kernel containing only this page. This kernel is now K_C , the current kernel under construction (at any instant there is just one of them). If all the pages were located in at least one kernel then go to point six.

⁵¹ (Whether absolute or with respect to the query)

⁵² In the case where we have the authority scores of the pages, or some other popularity score, we prefer in fact to base ourselves on them.

4. Find the relevant pages with respect to K_C (using the algorithm for calculating relevance scores) with

- $H = R^*$
- $A = R$
- $K = K_C$
- $T = R_X$

5. Let P_N be the page of R , not yet in K_C , which has the highest relevance score. If its relevance score is less than a fixed minimum score, return to point 3. (The current kernel is now complete). Otherwise insert it into K_C and go back to point four. It should be noted that it will not be necessary to reinitialize the hub and authority scores, it is preferable to keep the latest values calculated, thus the convergence ought to be very fast.

6. We now have a set of kernels (upstream homogeneous sub-queries) ready to be used as described in this document. When we want to calculate the relevance scores globally to the whole query we calculate an arithmetic average of the results for each of the kernels.

homogeneity equation

As a variant, instead of basing ourselves on the

$$\overline{x_S} = \prod_{P \in S} \overline{p_P}^{(-1)^{|P|}} \text{ as described}$$

example

method can be based on another homogeneity equation, such as for

$$x_S = \frac{\left| \bigcap_{P_i \in S} B(P_i) \right|}{\left| \bigcup_{P_i \in S} B(P_i) \right|} \text{ or else } x_S = \frac{\left| \bigcap_{P_i \in S} B(P_i) \right|}{\left| \bigcup_{P_i \in S} B(P_i) \right|} \cdot \left(\frac{\overline{\text{Min}_{P_i \in S} |B(P_i)|}}{\overline{\text{Max}_{P_i \in S} |B(P_i)|}} \right)$$

in which the ensemble cardinalities (represented between vertical bars) are replaced by the total of the hub scores of the pages in question⁵³.

Downstream processing:

Instead of searching for the good pages in relation to those of a kernel from among

⁵³ We can say that the cardinalities are replaced by "weighted cardinalities", the weights being the hub scores.

1 the pages that are cited in common with them it may be beneficial to perform the
2 same algorithms in the other direction, i.e. by searching among the pages which
3 cite the same pages as the kernel, or even to perform both and to calculate an
4 arithmetic average.

5 The downstream algorithms are identical to those upstream except that B is
6 replaced by F and F is replaced by B, and $-$ is interchanged with $+$ (for example R^{+}
7 is replaced by R^{-}).

8 The upstream and downstream methods may advantageously be integrated in the
9 following manner: after the upstream processing (possibly even after each
10 upstream iteration), with the candidate pages (R^{+}) having obtained a sufficient
11 relevance score, we associate downstream a set of extra pages ("artificial pages")
12 whose cardinality is dependent on said relevance score. Each artificial page is also
13 cited by (at least) one page of the query. The scores of these good pages (of R^{+})
14 found upstream⁵⁴ are thus given downstream an "advantage", and consequently the
15 scores of the pages (of R^{++}) cited as appropriate by these good pages are also
16 indirectly given an advantage.

17 And conversely, after the downstream processing (possibly even after each
18 downstream iteration), the same method is applied symmetrically upstream. Thus
19 the good pages of R^{+} are favored, as are indirectly the pages (of R^{+-}) which cite
20 them, as appropriate.

21 By not amalgamating the scores upstream (of the pages R^{+}) with the scores
22 downstream (pages R^{+}) it is possible to dissociate them in the calculations. In
23 particular, the influence of the scores obtained downstream can be decreased in the
24 upstream processing or vice versa.

25 Moreover, by virtue of this idea of "artificial pages", the present method may be
26 applied as a complement to the existing procedures of the prior art. Specifically,
27 once the scores have been obtained for each page, the respective numbers of citing
28 and cited pages can be modified artificially before applying these procedures.

29 It is possible to trek (known as "crawling") the web by following the links
30 (upstream and downstream) around the previously cited pages of the 7 sets,
31 exploiting the addition of the artificial pages to advantage the web pages linked to
32 the pages which are more relevant with respect to the query.

33 Insofar as the pages having the best scores are presumed to be relevant to the user
34 (and insofar as the relevance is transitive), the methods described here will be able

⁵⁴ Note that, advantageously, this is done without amalgamating the relevance scores upstream and downstream.

1 to be applied recursively thereto to discover yet other relevant pages. It is thus
2 possible to trek the web based on the user's query.

3 Figure 7 diagrammatically presents such a method: the search for relevant pages
4 can be applied recursively by extending the query with the "good pages found
5 upstream", "good pages found downstream", "good hub pages" and "good
6 authority pages" which in the figure are framed by rectangles. At each recursion,
7 the scores of the best pages found become slightly lower (because each time the
8 best pages found are added into the query) and the method stops when the scores
9 cease to be sufficient.

10 A system implementing the relative distillation method described hereinabove is
11 able to receive a search query composed of a set of URIs making it possible to
12 access information resources such as web pages and to provide in response the
13 URIs (or directly the pages) which are presumed to be the most relevant with
14 respect to said query.

15 The query being composed for example of the favorite links of the user and the
16 goal of the system being for example to monitor the web around these links and to
17 notify the user when new interesting pages appear therein, either by "Push"
18 technology at the initiative of a server, or by "Pull" technology at the initiative of
19 the user.

20 The user can of course provide the system directly with a set of URIs,
21 nevertheless, other means may also be offered to him to assist him in the
22 preparation and submission of a search query.

23 To trigger the execution of a search query from a hypertext link located in a page,
24 the user can use any one of the devices from among the following:

- 25 • A graphical object activatable for example by clicking (e.g. a button), is
26 presented close to certain hypertext links (URI) in a web page. Its activation
27 triggers the sending of a search query containing the URI in question.
- 28 • The system is furnished with a means able to toggle the page into a state
29 where each click on a link triggers the execution of a search query (containing this
30 link).
- 31 • A key of the keyboard, such as the "Ctrl" key, pressed while clicking (by a
32 means of pointing) serves to trigger the execution of a search query from the link
33 on which cursor of the pointing means is positioned.
- 34 • The right-hand mouse button (or equivalent) serves to trigger the execution
35 of a search query from the link on which the cursor of the mouse is positioned.

- 1 • Other analogous device.

2 Each of these devices can advantageously make it possible to execute said search
3 query in addition to (in parallel with) access to the page designated by the link in
4 question. The result of the search query will for example be displayed in a second
5 window (new instance of the browser) or else in a subwindow of the browser⁵⁵.

6 As a supplement to the link selected, other URIs may be added routinely into the
7 search query⁵⁶. They may in particular be:

- 8 • the links located in the page, in the region of the URI selected;
- 9 • the URIs previously selected by the user for this same query in the course
10 of his browsing⁵⁷;
- 11 • links explicitly envisaged and preferably determined by the designer of the
12 page to accompany the URI selected;
- 13 • the URIs that another user ("mentor" or referent) considers to be very
14 relevant with respect to the URI selected, the mentor being determined
15 automatically by the system, or specified by the user himself (chosen from a list of
16 "pals" that he has previously stored in the system), or else proposed by the page
17 designer (the user can also choose from a list of "experts" proposed by the page
18 designer).

19 Preparation of a query:

20 We shall now describe how the user can prepare a query composed of several links
21 that he gleans in the course of his browsing.

22

23 a) Displaying of the current query under preparation

24 Instead of triggering a search query directly, the user's action (as described above,
25 for example the act of clicking on a link with the right-hand button and choosing
26 the appropriate option) triggers the displaying of an accessory page in which:

- 27 • in addition to the link that the user wishes to select⁵⁸, other links, that he
28 has, as appropriate, previously selected for this same query, are presented;

29 — boxes to be ticked may be displayed in association with each link

⁵⁵ In a manner analogous to the subwindow existing today for favorites, this subwindow may be adjacent to the main subwindow in which the page containing the link that the user has clicked was displayed and in which the page accessed by the act of clicking on this link is subsequently displayed.

⁵⁶ Specifically, one of the essential advantages of the system is to be able to operate (find the relevant information resources) even if the search query is composed of a plurality of URIs.

⁵⁷ The new URIs found by the system are then highlighted in the result returned to the user (to distinguish them from the URIs which had already been returned in the same browsing).

⁵⁸ (As well as the links added routinely, as appropriate, as described hereinabove)

1 presented, in such a way that the user can in particular select those links that will
2 actually form the query;

3 • said accessory page is also furnished with an input means (such as a button)
4 making it possible to launch the search query.

5
6 Thus the user can prepare a query gradually, by selecting links one after the other⁵⁹
7 during his browsing⁶⁰ and thereafter send a query composed of several URIs.

8 Said accessory page may additionally contain drop-down graphical objects (such
9 as for example directories, records, folders, or similar metaphor) representing
10 queries under preparation other than the query in progress. The user can thus
11 choose the query or queries which will be enriched by the new link that he has just
12 selected.

13 Following the preparation of a query from a URI corresponding to a hypertext link
14 in a page (as described above), the already existing queries which, as appropriate,
15 contain this URI are optionally presented to him.

16 Advantageously, said accessory page may be composed of two parts. One of these
17 parts contains the elements described hereinabove (that is to say the elements of
18 the query under preparation). The other part presents the content of the page
19 designated by the link selected by the user.

20 For example, if the user clicks on a link while the page is in the state where all
21 clicks trigger the displaying of the current query under preparation (or with the
22 right-hand button of the mouse, etc.), the server returns said accessory page to it,
23 which thus comprises:

- 24 • in one part: the elements of the query under preparation
25 • and in the other part: the content of the page designated by the link clicked.

26 Thus, the use of the system represents an important advantage with respect to
27 conventional browsing around the web: the user receives not only the page
28 designated by the link that he has clicked (this is conventional web browsing), but
29 at the same time he benefits from the possibility of sending a query (containing
30 several URIs) to obtain yet other resources relevant in relation to this page.

31 As a variant, said accessory page is returned after fast (or even restricted⁶¹)
32 execution of the search query in the course of which the link clicked was added.

⁵⁹ (In one and the same page or in different pages)

⁶⁰ (During one and the same browsing or staggered over time)

⁶¹ In the case of a query regarding pages already crawled, the system can directly return the relevant URIs (or pages) already known and return the rest of the results later on.

1 The second page then directly contains a part of the result⁶². The user then receives
2 not only the page designated by the link that he has clicked, but in addition he
3 benefits directly from other resources relevant in relation to this page.

4 More advantageously still, said accessory page may be displayed in a subwindow⁶³
5 adjacent to the main subwindow of the browser. This adjacent subwindow opens in
6 response to the action of the user who desires the displaying of the query under
7 preparation (that is to say said accessory page).⁶⁴

8 The query under preparation can thus be displayed in parallel (asynchronously)
9 with the displaying of the page designated by the link clicked; the latter page being
10 displayed (independently) in the main subwindow.

11 The result of the search query can thereafter be presented in the same adjacent
12 subwindow.

13 As mentioned previously, a (partial) result may possibly be returned after partial or
14 restricted execution of the search query in progress, to which query the link clicked
15 was added. The adjacent subwindow then directly presents a fast search result
16 (which will possibly be supplemented subsequently).

17

18 b) Result of the execution of a search query

19 For each search query, the server can return the results directly (for example
20 returned from the HTTP query) or later on (for example by email).

21 The server returns the URIs (resulting from a query) in a page exhibiting the same
22 structure as said accessory page (or said query under preparation), namely:

23 • boxes to be ticked are associated with the links in such a way that the user
24 can select those links that he likes and delete those he does not like⁶⁵

25 – each URI⁶⁶ can thus be in at least one of the following states⁶⁷:
26 suggested (default state), accepted or deleted (the URIs that are in the deleted state
27 are not presented);

28 • the page is furnished with an input means (such as a button) making it

⁶² (For example in the form of a list of URIs or a set of vignettes representing these pages in miniature)

⁶³ (Analogous to the favorites subwindow of the current browsers)

⁶⁴ Note that, in parallel with the displaying of the query under preparation, the server can advantageously already begin to trek the web (crawling) – that is to say construct R^- , R^+ , R^{++} , R^+ , R^{+-} and R^{++} as already described – around the link selected.

⁶⁵ (That is to say ask the system to no longer suggest them)

⁶⁶ Optionally, the presentation of the result of a search query includes the content of the pages (that are pointed at by the resulting URIs) for example in miniaturized form (vignettes).

⁶⁷ Subsidiarily, an option to copy ("freeze") a page (locally or in a personal space on a server) may also be offered to the user. Each link can then be in one of the following states: suggested, accepted, deleted or frozen.

1 possible to relaunch the search query.

2 The page returned also presents the other queries (from the same user) in the form
3 of drop-down graphical objects, as already described. Their presentation may be
4 hierarchized according to their relevance with respect to the link clicked
5 (according to the relevance calculation methods described later).

6 The page returned presents means of control allowing the user to create new
7 queries and to delete existing queries. Of course, the user can cut and paste URIs
8 from existing queries or from any other resource. Also, when the result of a query
9 is returned by the server, the user can shift (hive off) the URIs received into other
10 queries. Each query is individually accessible by means of its own URI.

11

12 *Maintaining the spots*

13 Described hitherto are several methods that use the relative distillation procedure,
14 starting from a query (e.g. the given associated links of a spot) composed of a set
15 of URIs, to determine and store relevant URIs (e.g. the completed associated links
16 of a spot) with respect to this query, together with their relevance scores. These
17 stored results are obtained on the basis of counting links located in the resources of
18 the sets R^{+-} , R^{++} , R^{+-} , R^{+-} , R^{+-} , R^{+-} ⁶⁸ etc. which are themselves stored at least in
19 part. Now, these sets vary over time (and the links located in the resources
20 constituting these sets also vary). The stored data must therefore be kept up to date
21 and the calculations must be redone when the data that they take as input vary
22 significantly.

23 Moreover, it is desirable to disclose new relevant resources even before links
24 pointing to them appear on the web. A method making it possible to do so will
25 now be described.

26 For each query (for example for each spot),

- 27 - select a first set of resources having the largest relevance scores (such as
28 the largest hub scores) for said request,
- 29 - determine the *relevant regions* (that is to say the regions possessing links to
30 resources whose scores are high on average) of said first set of resources having
31 the largest relevance scores,
- 32 - monitor the new links which appear in said relevant regions and which
33 point to new resources (that is to say to resources that were not yet known to the
34 system),

⁶⁸ R^{+-} , R^{++} , and R^{+-} are in particular used to calculate the closeness of linked resources, and to filter, as described above, by taking the complement of this closeness as weighting for the counting of the links in question.

- 1 - select a second set of resources having a high relevance score (such as the
2 authority score) for said query,
3 - select the new resources which are the most similar to the resources of said
4 second set of resources and give the new resources selected a *time-dependent*
5 *authority score* (as described hereinbelow) as a function of their similarity to the
6 resources of said second set of resources.

7 The similarity of a resource with respect to other resources is determined by
8 comparing their contents. Described hereinbelow is the way to determine the
9 similarity as a function of the distribution of the words in the resources in question.

10

11 Time-dependent authority score:

12 Each new authority resource has a hypertext authority score (a_{ht}) and a similarity
13 authority score (a_s). Let τ be the ratio between

- 14 - the time remaining in order for the resource in question to no longer be
15 considered to be new

- 16 - and the total duration of newness (that is to say the total duration for which
17 a resource which has just been discovered by the system is considered to be new).

18 τ is therefore a number equal to 1 at the start of the life of a resource in the system,
19 and decreases linearly until it reaches 0 at the moment at which the resource in
20 question is said to be old.

21 Thus τ is used as a weighting to go gradually from a similarity score to a hypertext
22 score and the formula for the global score is $a = \tau a_s + \tau' a_{ht}$ (with $\tau' = 1 -$
23 τ).

24

25 As the distribution of the words of a new resource varies in principle less than the
26 hypertext links which point to it, a_s is considered to be constant while a_{ht} must be
27 updated over time. Thus the score a_s must be calculated at the moment at which
28 the new resource is discovered, and for all the queries for which it is in a relevant
29 region, until it becomes old (thus if a link to this resource appears in a relevant
30 region after it has become old, then its similarity with the resources of said second
31 set will not be determined).

32

33 Similarity:

34 An absolute distillation algorithm will be used to determine the score a_s of each
35 new resource.

36 The known method of absolute distillation over a set of nodes connected by links
37 (thus forming an oriented graph) comprises the following steps:

- 1 1 - allocate each node a hub score equal to 1 and an authority score,
- 2 2 - for each node calculate its authority score by adding up the hub scores of
- 3 the nodes which point to it, then normalize the authority scores in such a way that
- 4 their total is equal to 1,
- 5 3 - for each node calculate its hub score by adding up the authority scores of
- 6 the nodes to which it points, then normalize the hub scores in such a way that their
- 7 total is equal to 1,
- 8 4 - reiterate by restarting from step 2 until the algorithm converges, that is to
- 9 say until the scores are no longer significantly different with respect to the
- 10 previous step.

11 In addition, here the links are weighted by the similarities of the resources in
12 question with respect to the distribution of their words. Steps 2 and 3 are replaced
13 by the following:

- 14 2' - for each node calculate its authority score by adding up the hub scores of
- 15 the nodes which point to it, multiplied by the weight of the respective links, then
- 16 normalize the authority scores in such a way that their total is equal to 1,
- 17 3' - for each node calculate its hub score by adding up the authority scores of
- 18 the nodes to which it points, multiplied by the weights of the respective links, then
- 19 normalize the hub scores in such a way that their total is equal to 1.

20 The weight of the similarity link between two resources is equal to the scalar
21 product of their distributions of words (that is to say to the sum, for each word
22 located in the two resources, of the product of the frequencies of this word in these
23 resources; the resulting sum is a number between zero – case where there is no
24 word in common – and 1 – case where the two resources have the same content)
25 after having removed the nonpertinent words ("stop words").

26 It should be noted that the similarity links thus obtained are bidirectional.

27 Thus, the absolute distillation can thus be performed over the set of resources
28 comprising:

- 29 - the new resource discovered,
- 30 - and said second set of resources having high relevance scores,
- 31 to determine the scores a_s of this new resource discovered.

32 The methods described above also make it possible to select, from among a set of
33 extra resources, a resource which is the most relevant with respect to a starting
34 resource.

- 1 Accordingly, the following three steps are implemented:
- 2 (a) selection from the web of resources that are most similar to the starting
- 3 resource (typically a private resource), by one of the procedures of the invention,
- 4 (b) selection from the web of resources that are the most relevant with respect
- 5 to the resources selected in step (a), and
- 6 (c) selection of extra resources (typically of private resources again) that are
- 7 the most similar to the most relevant resources selected in step (b).
- 8
- 9 Such a method makes it possible in particular to dynamically generate the content
- 10 of web pages published as a function of context.